

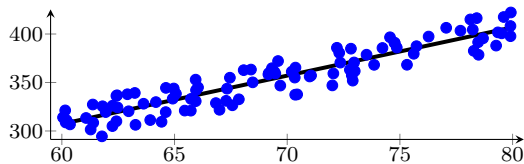
4.1: Visualizing Variability with a Scatterplot

Definition.

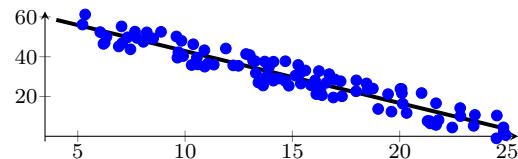
A **scatterplot** is used when examining the relationship between two *numerical* variables where each point represents an observation. With scatterplots, we examine the

- **trend:** general tendency of the scatter plot going from left to right
- **strength:** strong associations have little vertical variation
- **shape:** is the scatterplot linear or nonlinear?

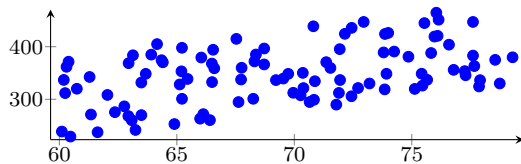
Positive trend



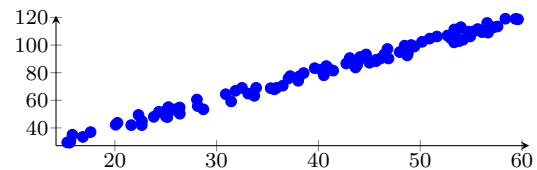
Negative trend



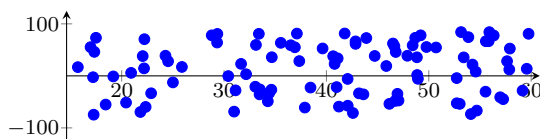
Weak association



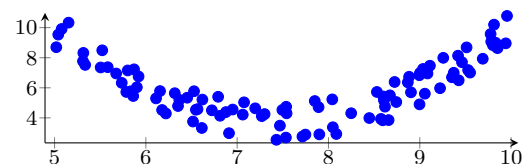
Strong association



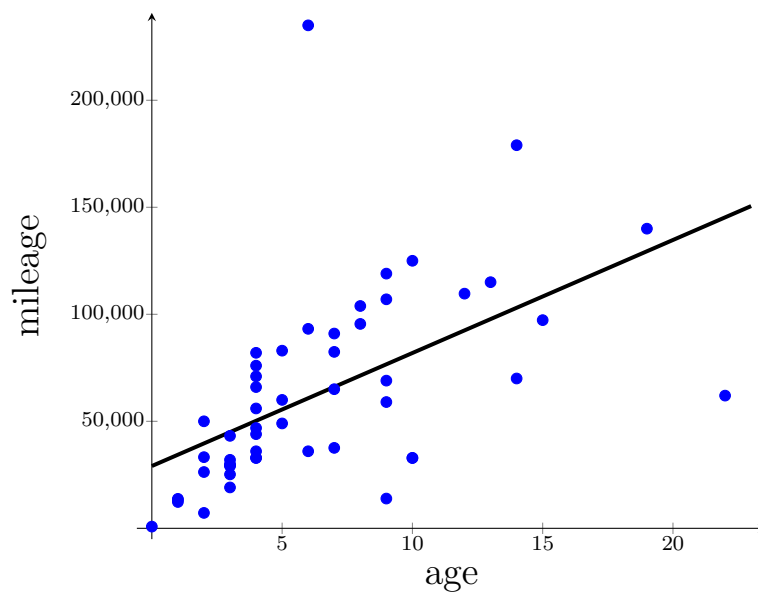
No trend



Quadratic/nonlinear shape



Example. The scatterplot below shows the age and corresponding mileage for a sample of used cars.



What is the association between the variables? Identify the trend, its shape, and how strong the relationship is.

4.2: Measuring Strength of Association with Correlation

Definition.

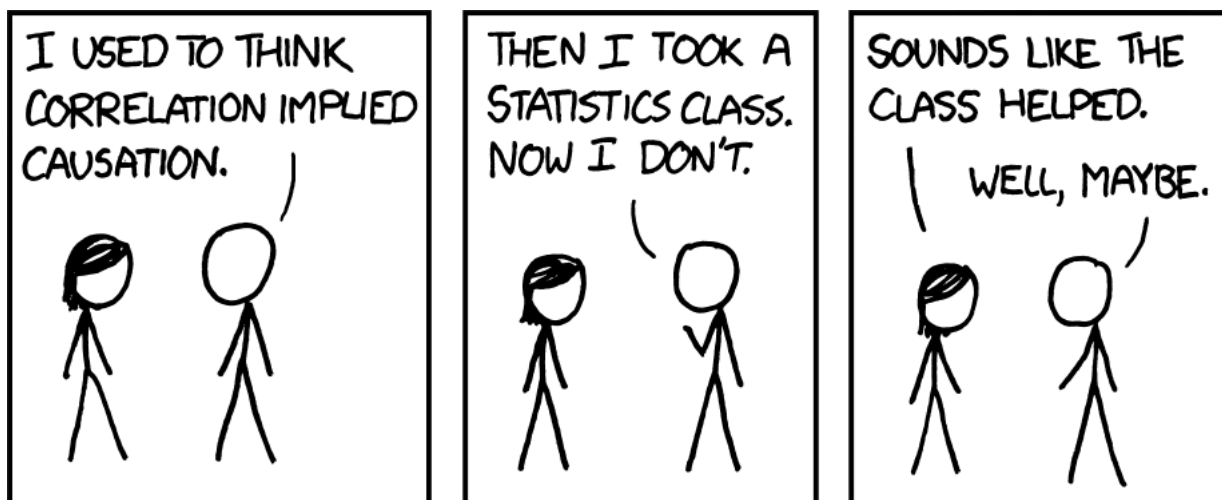
The **correlation coefficient** is a number that measures the strength of the linear association between two numerical variables. The correlation coefficient is between -1 and 1 :

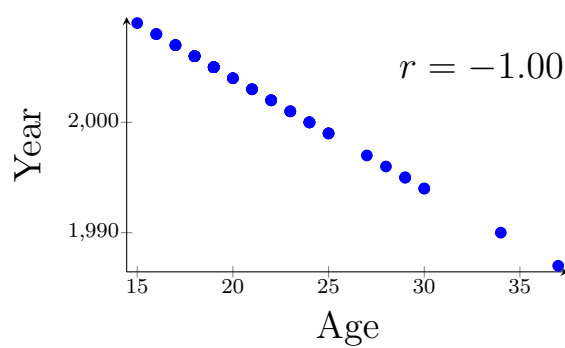
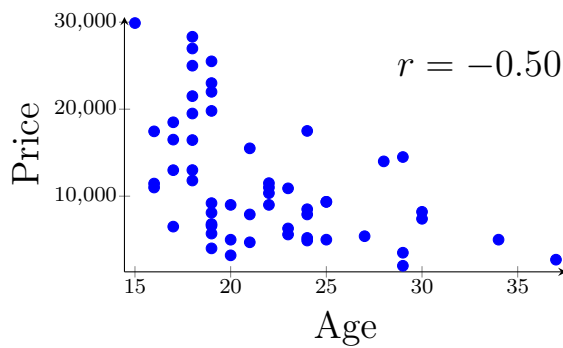
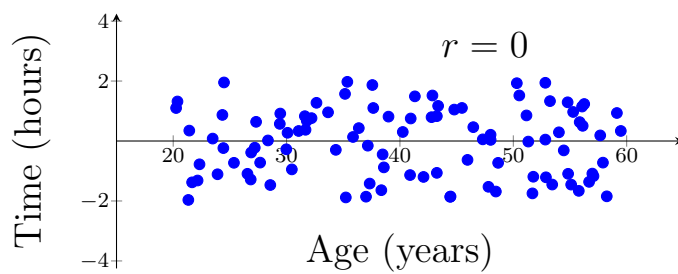
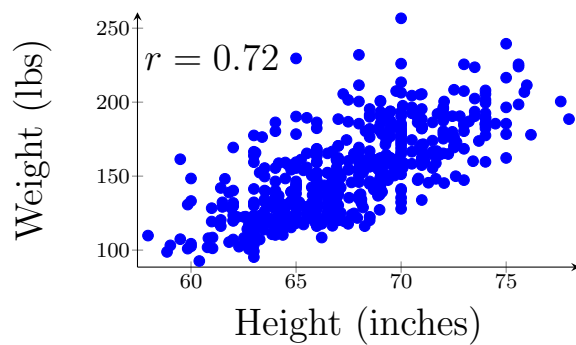
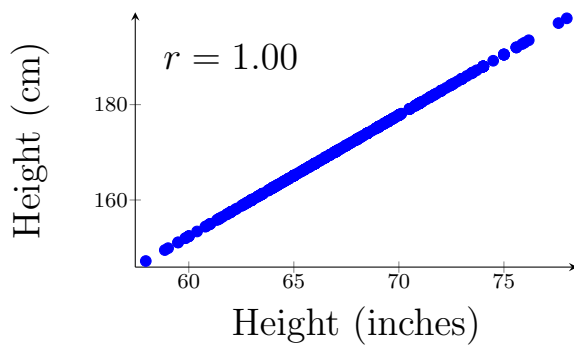
$1 \rightarrow$	Strong association and positive trend
$0 \rightarrow$	Weak or no association
$-1 \rightarrow$	Strong association and negative trend

The correlation coefficient only makes sense if the trend is linear and both variables are numerical!!

Note: *Correlation does not mean causation!*

Take a few minutes to look at the graphs at this link: [Spurious correlations](#)





Definition.

The formula for the correlation coefficient between two variables x and y is

$$r = \frac{\sum z_x z_y}{n - 1}$$

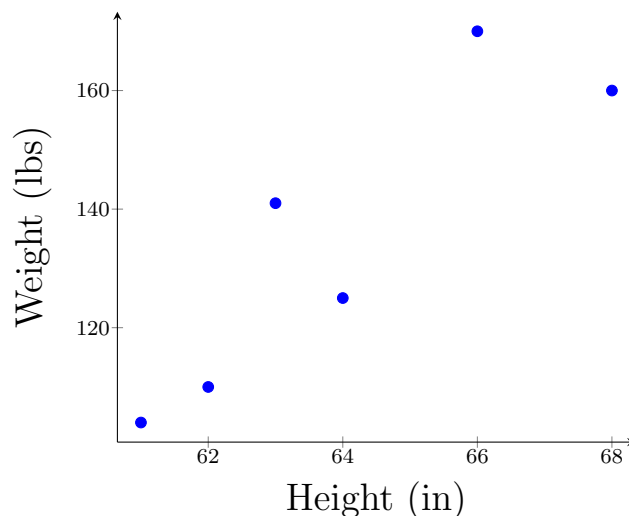
where z_x and z_y are the z scores for each entry in the x and y lists.

$$z_x = \frac{x - \bar{x}}{s_x} \quad z_y = \frac{y - \bar{y}}{s_y}$$

Example. Below are the heights and weights of six women:

Heights	61	62	63	64	66	68
Weights	104	110	141	125	170	160

Compute the correlation coefficient by hand. Then, graph the scatterplot and compute the correlation coefficient using StatCrunch.



Understanding the Correlation Coefficient:

- Changing the order of the variables does not change r
- Adding a constant or multiplying by a positive constant does not affect r
- The correlation coefficient is unitless
- None of this makes any sense if the relationship between the variables is not linear!